

# Learning with $\ell^0$ -Graph: $\ell^0$ -Induced Sparse Subspace Clustering

## Abstract

Sparse subspace clustering methods, such as Sparse Subspace Clustering (SSC) [8] and  $\ell^1$ -graph [26, 4], are effective in partitioning the data that lie in a union of subspaces. Most of those methods use  $\ell^1$ -norm or  $\ell^2$ -norm with thresholding to impose the sparsity of the constructed sparse similarity graph, and certain assumptions, e.g. independence or disjointness, on the subspaces are required to obtain the subspace-sparse representation, which is the key to their success. Such assumptions are not guaranteed to hold in practice and they limit the application of sparse subspace clustering on subspaces with general location. In this paper, we propose a new sparse subspace clustering method named  $\ell^0$ -graph. In contrast to the required assumptions on subspaces for most existing sparse subspace clustering methods, it is proved that subspace-sparse representation can be obtained by  $\ell^0$ -graph for arbitrary distinct underlying subspaces almost surely under the mild i.i.d. assumption on the data generation. We develop a proximal method to obtain the sub-optimal solution to the optimization problem of  $\ell^0$ -graph with proved guarantee of convergence. Moreover, we propose a regularized  $\ell^0$ -graph that encourages nearby data to have similar neighbors so that the similarity graph is more aligned within each cluster and the graph connectivity issue is alleviated. Extensive experimental results on various data sets demonstrate the superiority of  $\ell^0$ -graph compared to other competing clustering methods, as well as the effectiveness of regularized  $\ell^0$ -graph.

## 1. Introduction

Clustering is a common unsupervised data analysis method which partitions data into a set of self-similar clusters. High dimensionality of data often imposes difficulty on clustering. For example, model-based clustering methods, such as Gaussian Mixture Model (GMM) that models the data by a mixture of parametric distributions, suffer from the curse of dimensionality when fitting a statistical model to the data [9].

Based on the observation that high dimensional data often lie in a set of low-dimensional subspaces in many practical scenarios, subspace clustering algorithms [24] aim to

partition the data such that data belonging to the same subspace are identified as one cluster. Among various subspace clustering algorithms, the ones that employ sparsity prior, such as Sparse Subspace Clustering (SSC) [8] and  $\ell^1$ -graph [26, 4], have been proven to be effective in separating the data in accordance with the subspaces that the data lie in under certain assumptions.

Sparse subspace clustering methods construct sparse similarity graph by sparse representation of the data, where the vertices represent the data, and an edge is between two vertices whenever one participates the sparse representation of the other. Thanks to the subspace-sparse representation, the nonzero elements in the sparse representation of each datum in a subspace correspond to the data points in the same subspace, so that vertices corresponding to different subspaces are disconnected in the sparse similarity graph, leading to their compelling performance with spectral clustering [18] applied on such graph.

[8] proves that when the subspaces are independent or disjoint, then subspace-sparse representations can be obtained by solving the canonical sparse coding problem using data as the dictionary under certain conditions on the rank, or singular value of the data matrix and the principle angle between the subspaces respectively. Under the independence assumption on the subspaces, low rank representation [13, 12] is also proposed to recover the subspace structures. Relaxing the assumptions on the subspaces to allowing overlapping subspaces, the Greedy Subspace Clustering [19] and the Low-Rank Sparse Subspace Clustering [25] achieve subspace-sparse representation with high probability. However, their results rely on the semi-random model which assumes the data in each subspace are generated i.i.d. uniformly on the unit sphere in that subspace as well as certain additional conditions on the size and dimensionality of the data. In addition, the geometric analysis in [22] also adopts the semi-random model and it handles overlapping subspaces.

To avoid the non-convex optimization problem incurred by  $\ell^0$ -norm, most of the sparse subspace clustering or sparse graph based clustering methods use  $\ell^1$ -norm [26, 4, 7, 8, 28] or  $\ell^2$ -norm with thresholding [20] to impose the sparsity on the constructed similarity graph. In addition,  $\ell^1$ -norm has been widely used as a convex relaxation of  $\ell^0$ -norm for efficient sparse coding algorithms [11, 14, 15].

On the other hand, sparse representation methods such as [16] that directly optimize objective function involving  $\ell^0$ -norm demonstrate compelling performance compared to its  $\ell^1$ -norm counterpart. It remains an interesting question whether sparse subspace clustering equipped with  $\ell^0$ -norm, which is the origination of the sparsity that counts the number of nonzero elements, has advantage in obtaining the subspace-sparse representation. In this paper, we propose  $\ell^0$ -graph which employs  $\ell^0$ -norm to enforce the sparsity of the similarity graph. This paper offers three contributions:

**Theoretical Results on  $\ell^0$ -Induced Almost Surely Subspace-Sparse Representation** We present the theory of the  $\ell^0$ -induced sparse subspace clustering by  $\ell^0$ -graph, which shows that  $\ell^0$ -graph renders subspace-sparse representation almost surely under minimum assumptions on the underlying subspaces the data lie in, i.e. subspaces are distinct. To the best of our knowledge, this is the mildest assumption on the subspaces compared to most existing sparse subspace clustering methods. Furthermore, our theory assumes that the data in each subspace are generated i.i.d. from arbitrary continuous distribution supported on that subspace, which is milder than the assumption of semi-random model in [19] and [25] that assume the data are i.i.d. uniformly distributed on the unit sphere in each subspace.

**Efficient Optimization** The optimization problem of  $\ell^0$ -graph is NP-hard and it is impractical to pursue the global optimal solution. Instead, we develop an efficient proximal method to obtain a sub-optimal solution with convergence guarantee.

**Regularized  $\ell^0$ -Graph** In order to obtain a sparse similarity graph where neighboring data have similar neighbors so as to encourage the graph connectivity within each cluster, we propose Regularized  $\ell^0$ -graph that incorporates a regularization term into the objective of  $\ell^0$ -graph. Moreover, we have implemented both  $\ell^0$ -graph and regularized  $\ell^0$ -graph in CUDA C programming language for significant speedup by parallel computing.

Note that SSC-OMP [6] adopts Orthogonal Matching Pursuit (OMP) [23] to choose neighbors for each datum in the sparse similarity graph, which can be interpreted as approximately solving a  $\ell^0$  problem. However, SSC-OMP does not present the theoretical properties of the  $\ell^0$ -induced sparse subspace clustering, and the experimental results show the significant performance advantage of  $\ell^0$ -graph over the OMP-graph. OMP-graph solves the  $\ell^0$  problem of  $\ell^0$ -graph by OMP, so that it is equivalent to SSC-OMP for clustering. Although our optimization algorithm only obtains a sub-optimal solution to the objective of  $\ell^0$ -graph,

we give theory about  $\ell^0$ -induced subspace structures and extensive experimental results show the effectiveness of our model.

The remaining parts of the paper are organized as follows. The representative subspace subspace clustering methods, SSC and  $\ell^1$ -graph, are introduced in the next subsection, and then the detailed formulation of  $\ell^0$ -graph and regularized  $\ell^0$ -graph is illustrated. We then show the clustering performance of the proposed models, and conclude the paper. We use bold letters for matrices and vectors, and regular lower letter for scalars throughout this paper. The bold letter with superscript indicates the corresponding column of a matrix, and the bold letter with subscript indicates the corresponding element of a matrix or vector.  $\|\cdot\|_F$  and  $\|\cdot\|_p$  denote the Frobenius norm and the  $\ell^p$ -norm, and  $\text{diag}(\cdot)$  indicates the diagonal elements of a matrix.

## 1.1. Sparse Subspace Clustering and $\ell^1$ -Graph

Sparse coding methods represent an input signal by a linear combination of only a few atoms of a dictionary, and the sparse coefficients are named sparse code. Sparse coding has been broadly applied in machine learning and signal processing, and sparse code is extensively used as a discriminative and robust feature representation [27, 5, 29, 28]

SSC [8] and  $\ell^1$ -graph [26, 4] employ sparse representation of the data to construct the sparse similarity graph. With the data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  where  $n$  is the size of the data and  $d$  is the dimensionality, SSC and  $\ell^1$ -graph solves the following sparse coding problem:

$$\min_{\alpha} \|\alpha\|_1 \quad \text{s.t. } \mathbf{X} = \mathbf{X}\alpha, \quad \text{diag}(\alpha) = \mathbf{0} \quad (1)$$

Both SSC and  $\ell^1$ -graph construct a sparse similarity graph  $G = (\mathbf{X}, \mathbf{W})$  where the data  $\mathbf{X}$  are represented as vertices,  $\mathbf{W}$  is the graph weight matrix of size  $n \times n$  and  $\mathbf{W}_{ij}$  indicates the similarity between  $\mathbf{x}_i$  and  $\mathbf{x}_j$ ,  $\mathbf{W}$  is set by the sparse codes  $\alpha$  as below:

$$\mathbf{W}_{ij} = (|\alpha_{ij}| + |\alpha_{ji}|)/2 \quad 1 \leq i, j \leq n \quad (2)$$

Furthermore, suppose the underlying subspaces that the data lie in are independent or disjoint, SSC [8] proves that the optimal solution to (1) is the subspace-sparse representation under several additional conditions. *The sparse representation  $\alpha$  is called subspace-sparse representation if the nonzero elements of  $\alpha^i$ , namely the sparse representation of the datum  $\mathbf{x}_i$ , correspond to the data points in the same subspace as  $\mathbf{x}_i$ .* Therefore, vertices corresponding to different subspaces are disconnected in the sparse similarity graph. With the subsequent spectral clustering [18] applied on such sparse similarity graph, compelling clustering performance is achieved.

Allowing some tolerance for inexact representation, the literature often turns to solve the following problem for SSC and  $\ell^1$ -graph:

$$\min_{\alpha} \|\alpha\|_1 \quad \text{s.t. } \|\mathbf{X} - \mathbf{X}\alpha\|_F \leq \delta, \quad \text{diag}(\alpha) = \mathbf{0}$$

which is equivalent to the following problem

$$\min_{\alpha} \|\mathbf{X} - \mathbf{X}\alpha\|_F^2 + \lambda_{\ell^1} \|\alpha\|_1 \quad s.t. \quad \text{diag}(\alpha) = \mathbf{0} \quad (3)$$

where  $\lambda_{\ell^1} > 0$  is a weighting parameter for the  $\ell^1$  term.

## 2. $\ell^0$ -Induced Sparse Subspace Clustering

In this paper, we investigate  $\ell^0$ -induced sparse subspace clustering method, which solves the following  $\ell^0$  problem:

$$\min_{\alpha} \|\alpha\|_0 \quad s.t. \quad \mathbf{X} = \mathbf{X}\alpha, \quad \text{diag}(\alpha) = \mathbf{0} \quad (4)$$

We then give the theorem about  $\ell^0$ -induced almost surely subspace-sparse representation, and the proof is presented in the supplementary document for this paper.

**Theorem 1.** ( *$\ell^0$ -Induced Almost Surely Subspace-Sparse Representation*) Suppose the data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  lie in a union of  $K$  distinct subspaces  $\{\mathcal{S}_k\}_{k=1}^K$  of dimensions  $\{d_k\}_{k=1}^K$ , i.e.  $\mathcal{S}_k \neq \mathcal{S}_{k'}$  for  $k \neq k'$ . Let  $\mathbf{X}^{(k)} \in \mathbb{R}^{d \times n_k}$  denotes the data that belong to subspace  $\mathcal{S}_k$ , and  $\sum_{k=1}^K n_k = n$ . When  $n_k \geq d_k + 1$ , if the data belonging to each subspace are generated i.i.d. from some unknown distribution supported on that subspace, then with probability 1, the optimal solution to (4), denoted by  $\alpha^*$ , is a subspace-sparse representation, i.e. nonzero elements in  $\alpha^{*i}$  corresponds to the data that lie in the same subspace as  $\mathbf{x}_i$ .

Based on the above theorem, we propose  $\ell^0$ -graph that solves (4) and uses the sparse representation to build the sparse similarity graph for clustering. According to Theorem 1,  $\ell^0$ -induced sparse subspace clustering method (4) obtains the subspace-sparse representation almost surely under minimum assumption on the subspaces, i.e. it only requires that the subspaces be distinct. To the best of our knowledge, this is the mildest assumption on the subspaces for most existing sparse subspace clustering methods. Moreover, the only assumption on the data generation is that the data in each subspace are i.i.d. random samples from arbitrary continuous distributions supported on that subspace. In the light of assumed data distribution, such assumption on the data generation is much milder than the assumption of the semi-random model in ([19, 25, 22]) (note that the data can always be normalized to have unit norm and reside on the unit sphere). Table 1 summarizes different assumptions on the subspaces and random data generation for different sparse subspace clustering methods. It can be seen that  $\ell^0$ -graph has mildest assumption on both subspaces and the random data generation.

## 3. Optimization of $\ell^0$ -Graph

We introduce the optimization algorithm for  $\ell^0$ -graph in this section. Similar to the case of SSC and  $\ell^1$ -graph, by

allowing tolerance for inexact representation, we turn to optimize the following  $\ell^0$  problem

$$\min_{\alpha} L(\alpha) = \|\mathbf{X} - \mathbf{X}\alpha\|_F^2 + \lambda \|\alpha\|_0 \quad s.t. \quad \text{diag} = \mathbf{0} \quad (5)$$

Problem (5) is NP-hard, and it is impractical to seek for its global optimal solution. The literature extensively resorts to approximate algorithms, such as Orthogonal Matching Pursuit [23], or that uses surrogate functions [10], for  $\ell^0$  problems. Inspired by recent advances in solving non-convex optimization problems by proximal linearized method [3] and the application of this method to  $\ell^0$ -norm based dictionary learning [2], we propose an iterative proximal method to optimize (5) and obtain a sub-optimal solution with proved convergence guarantee. In the following text, the superscript with bracket indicates the iteration number of the proposed proximal method.

In  $t$ -th iteration of our proximal method for  $t \geq 1$ , gradient descent is performed on the squared loss term of (5), i.e.  $Q(\alpha) = \|\mathbf{X} - \mathbf{X}\alpha\|_F^2$ , to obtain

$$\tilde{\alpha}^{(t)} = \alpha^{(t-1)} - \frac{2}{\tau s} (\mathbf{X}^\top \mathbf{X} \alpha^{(t-1)} - \mathbf{X}^\top \mathbf{X}) \quad (6)$$

where  $\tau$  is any constant that is greater than 1, and  $s$  is the Lipschitz constant for the gradient of function  $Q(\cdot)$ , namely

$$\|\nabla Q(\mathbf{Y}) - \nabla Q(\mathbf{Z})\|_F \leq s \|\mathbf{Y} - \mathbf{Z}\|_F, \quad \forall \mathbf{Y}, \mathbf{Z} \in \mathbb{R}^{n \times n} \quad (7)$$

Then  $\alpha^{(t)}$  is the solution to the following  $\ell^0$  regularized problem:

$$\begin{aligned} \alpha^{(t)} &= \arg \min_{\mathbf{v} \in \mathbb{R}^{n \times n}} \frac{\tau s}{2} \|\mathbf{v} - \tilde{\alpha}^{(t)}\|_F^2 + \lambda \|\mathbf{v}\|_0 \\ s.t. \quad &\text{diag}(\mathbf{v}) = \mathbf{0} \end{aligned} \quad (8)$$

It can be verified that (8) has closed-form solution, i.e.

$$\alpha_{ij}^{(t)} = \begin{cases} 0 & : |\tilde{\alpha}_{ij}^{(t)}| < \sqrt{\frac{2\lambda}{\tau s}} \text{ or } i = j \\ \tilde{\alpha}_{ij}^{(t)} & : \text{otherwise} \end{cases} \quad (9)$$

for  $1 \leq i, j \leq n$ . The iterations start from  $t = 1$  and continue until the sequence  $\{L(\alpha^{(t)})\}$  converges or maximum iteration number is achieved. We initialize  $\alpha$  as  $\alpha^{(0)} = \alpha_{\ell^1}$  and  $\alpha_{\ell^1}$  is the sparse codes generated by SSC or  $\ell^1$ -graph via solving (3) with some proper weighting parameter  $\lambda_{\ell^1}$ . In all the experimental results of this paper, we empirically set  $\lambda_{\ell^1} = 0.1$  when initializing  $\ell^0$ -graph.

The data clustering algorithm by  $\ell^0$ -graph is described in Algorithm 1. Also, the following theorem shows that each iteration of the proposed proximal method decreases the value of the objective function  $L(\cdot)$  in (5), therefore, our proximal method always converges.

**Theorem 2.** Let  $s = 2\sigma_{\max}(\mathbf{X}^\top \mathbf{X})$  where  $\sigma_{\max}(\cdot)$  indicates the largest eigenvalue of a matrix, then the sequence  $\{L(\alpha^{(t)})\}$  generated by the proximal method with (6) and (9) decreases, and the following inequality holds for  $t \geq 1$ :

$$L(\alpha^{(t)}) \leq L(\alpha^{(t-1)}) - \frac{(\tau - 1)s}{2} \|\alpha^{(t)} - \alpha^{(t-1)}\|_F^2 \quad (10)$$

Table 1. Assumptions on the subspaces and random data generation (for randomized part of the algorithm) for different sparse subspace clustering methods. Note that  $S_1 < S_2 < S_3 < S_4$ ,  $D_1 < D_2$ , and the assumption on the right hand side of  $<$  is milder than that on the left hand side.

Assumption on Subspaces	Explanation
$S_1$ :Independent Subspaces ([13, 12])	$\text{Dim}[S_1 \otimes S_2 \dots S_K] = \sum_k \text{Dim}[S_k]$
$S_2$ :Disjoint Subspaces ([8])	$S_k \cap S_{k'} = \mathbf{0}$ for $k \neq k'$
$S_3$ :Overlapping Subspaces ([19, 25, 22])	$\text{Dim}[S_k \cap S_{k'}] < \min\{\text{Dim}[S_k], \text{Dim}[S_{k'}]\}$ for $k \neq k'$
$S_4$ :Distinct Subspaces ( $\ell^0$ -Graph)	$S_k \neq S_{k'}$ for $k \neq k'$
Assumption on Random Data Generation	Explanation
$D_1$ :Semi-Random Model ([19, 25, 22])	The data in each subspace are generated i.i.d. uniformly on the unit sphere in that subspace.
$D_2$ :IID ( $\ell^0$ -Graph)	The data in each subspace are generated i.i.d. from arbitrary continuous distribution supported on that subspace.

And it follows that the sequence  $\{L(\alpha^{(t)})\}$  converges.

Furthermore, we show that if the sequence  $\{\alpha^{(t)}\}$  generated by the proposed proximal method is bounded, then it is a Cauchy sequence and it converges to a critical point of the objective function  $L$  in (5).

**Theorem 3.** Suppose that the sequence  $\{\alpha^{(t)}\}$  generated by the proximal method with (6) and (9) is bounded, then 1)  $\sum_{t=1}^{\infty} \|\alpha^{(t)} - \alpha^{(t-1)}\|_F < \infty$  2)  $\{\alpha^{(t)}\}$  converges to a critical point<sup>1</sup> of the function  $L(\cdot)$  in (5).

*Sketch of the Proof.* [3] shows that the  $\ell^0$ -norm function  $\|\cdot\|_0$  is a semi-algebraic function. The conclusions of this theorem directly follows from Theorem 1 in [3].  $\square$

The detailed proofs of Theorem 2 and Theorem 3 are included in the supplementary document.

#### 4. Regularized $\ell^0$ -Graph

While the subspace-sparse representation separates the data belonging to different subspaces in the constructed sparse similarity graph, it is not guaranteed that the data points in the same subspace form a connected component. This is the well known graph connectivity issue in the sparse subspace clustering literature [8, 17] which is the only gap that prevents a sparse similarity graph with subspace-sparse representation from forming the perfect clustering result, i.e. the data belonging to each subspace form a single connected component in the sparse similarity graph. SSC [8] suggests alleviating the graph connectivity issue by promoting common neighbors across the data in each subspace. In this section we propose Regularized  $\ell^0$ -Graph by adding a regularization term to (5) which employs  $\ell^0$ -distance between the sparse representation of the data so as to impose the sparsity of the representation and encourage common neighbors for nearby data simultaneously. Regularized  $\ell^0$ -graph solves the following problem

$$\min_{\alpha} \|\mathbf{X} - \mathbf{X}\alpha\|_F^2 + \gamma R_S(\alpha) \quad (11)$$

<sup>1</sup> $x$  is a critical point of function  $f$  if  $0 \in \partial f(x)$ , where  $\partial f(x)$  is the limiting-subdifferential of  $f$  at  $x$ . Please refer to more detailed definition in [3].

#### Algorithm 1 Data Clustering by $\ell^0$ -Graph

##### Input:

- The data set  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ , the number of clusters  $c$ , the parameter  $\lambda$  for  $\ell^0$ -graph,  $\lambda_{\ell^1}$  for the initialization of the  $\ell^0$ -graph, maximum iteration number  $M$ , stopping threshold  $\varepsilon$
- 1:  $t = 1$ , initialize the coefficient matrix as  $\alpha^{(0)} = \alpha_{\ell^1}$ ,  $s = 2\sigma_{\max}(\mathbf{X}^\top \mathbf{X})$ .
  - 2: **while**  $t \leq M$  **do**
  - 3:   Obtain  $\alpha^{(t)}$  from  $\alpha^{(t-1)}$  by (6) and (9)
  - 4:   **if**  $|L(\alpha^{(t)}) - L(\alpha^{(t-1)})| < \varepsilon$  **then**
  - 5:     **break**
  - 6:   **else**
  - 7:      $t = t + 1$ .
  - 8:   **end if**
  - 9: **end while**
  - 10: Obtain the sub-optimal coefficient matrix  $\alpha^*$  when the above iterations converge or maximum iteration number is achieved.
  - 11: Build the sparse similarity matrix by symmetrizing  $\alpha^*$ :  $\mathbf{W}^* = \frac{|\alpha^*| + |\alpha^*|^\top}{2}$ , compute the corresponding normalized graph Laplacian  $\mathbf{L}^* = (\mathbf{D}^*)^{-\frac{1}{2}}(\mathbf{D}^* - \mathbf{W}^*)(\mathbf{D}^*)^{-\frac{1}{2}}$ , where  $\mathbf{D}^*$  is a diagonal matrix with  $\mathbf{D}_{ii}^* = \sum_{j=1}^n \mathbf{W}_{ij}^*$
  - 12: Construct the matrix  $\mathbf{v} = [\mathbf{v}_1, \dots, \mathbf{v}_c] \in \mathbb{R}^{n \times c}$ , where  $\{\mathbf{v}_1, \dots, \mathbf{v}_c\}$  are the  $c$  eigenvectors of  $\mathbf{L}^*$  corresponding to its  $c$  smallest eigenvalues. Treat each row of  $\mathbf{v}$  as a data point in  $\mathbb{R}^c$ , and run K-means clustering method to obtain the cluster labels for all the rows of  $\mathbf{v}$ .

**Output:** The cluster label of  $\mathbf{x}_i$  is set as the cluster label of the  $i$ -th row of  $\mathbf{v}$ ,  $1 \leq i \leq n$ .

where  $R_S(\alpha) = \sum_{i,j=1}^n \mathbf{S}_{ij} \|\alpha^i - \alpha^j\|_0$  is the regularization term,  $\mathbf{S}$  is the adjacency matrix of the KNN graph and  $\mathbf{S}_{ij} = 1$  if and only if  $\mathbf{x}_i$  is among the  $K$  nearest neighbors of  $\mathbf{x}_j$  in the sense of Euclidean distance. It should be emphasized that such KNN graph is a widely used strategy to identify nearby data for graph regularization in sparse



coding [30, 28].  $\gamma > 0$  is the weighting parameter for the regularization term. Since the co-located elements of two sparse codes  $\alpha^i$  and  $\alpha^j$  are not exactly the same in most cases, their  $\ell^0$ -distance  $\|\alpha^i - \alpha^j\|_0$  is almost always the sum of their difference in support and the number of their co-located nonzero elements, and the support of a vector is defined to be the indices of its nonzero elements. Therefore, the regularization term  $R_{\tilde{\mathbf{S}}}(\alpha)$  encourages both sparsity and common neighbors across nearby data.

We use coordinate descent to optimize (11) with respect to  $\alpha^i$  in each step of the coordinate descent, with all the other sparse codes  $\{\alpha^j\}_{j \neq i}$  fixed. The optimization problem for  $\alpha^i$  in each step is presented below:

$$\min_{\alpha^i} F(\alpha^i) = \|\mathbf{x}_i - \mathbf{X}\alpha^i\|_2^2 + \gamma R_{\tilde{\mathbf{S}}}(\alpha^i) \quad (12)$$

where  $R_{\tilde{\mathbf{S}}}(\alpha^i) = \sum_{j=1}^n \tilde{\mathbf{S}}_{ij} \|\alpha^i - \alpha^j\|_0$ , where  $\tilde{\mathbf{S}} = \mathbf{S} + \mathbf{S}^\top$ .

(12) can also be optimized by the proximal method in a similar manner to  $\ell^0$ -graph. In  $t$ -th ( $t \geq 1$ ) iteration of our proximal method for the problem (12), gradient descent on the squared loss term of the objective function of (12) is performed by (13):

$$\tilde{\alpha}^{i(t)} = \alpha^{i(t-1)} - \frac{2}{\tau s} (\mathbf{X}^\top \mathbf{X} \alpha^{i(t-1)} - \mathbf{X}^\top \mathbf{x}_i) \quad (13)$$

where  $\tau$  and  $s$  are the same as that in (6). Then  $\alpha^{i(t)}$  is obtained as the solution to the following  $\ell^0$  regularized problem:

$$\alpha^{i(t)} = \arg \min_{\mathbf{v} \in \mathbb{R}^n, \mathbf{v}_i = 0} \frac{\tau s}{2} \|\mathbf{v} - \tilde{\alpha}^{i(t)}\|_2^2 + \gamma R_{\tilde{\mathbf{S}}}(\mathbf{v}) \quad (14)$$

Proposition 1 below shows the closed form solution to the subproblem (14):

**Proposition 1.** Define  $F_k(v) = \frac{\tau s}{2} \|v - \tilde{\alpha}_{ki}^{(t)}\|_2^2 + \gamma R_{\tilde{\mathbf{S}}}(v)$  for  $v \in \mathbb{R}$  and  $R_{\tilde{\mathbf{S}}}(v) \triangleq \sum_{j=1}^n \tilde{\mathbf{S}}_{ij} \|v - \alpha_{kj}\|_0$ . Let  $\mathbf{v}^*$  be the optimal solution to (14), then the  $k$ -th element of  $\mathbf{v}^*$  is

$$\mathbf{v}_k^* = \begin{cases} \arg \min_{v \in \{\tilde{\alpha}_{ki}^{(t)}\} \cup \{\alpha_{kj}\}_{j: \tilde{\mathbf{S}}_{ij} \neq 0}} F_k(v) & : k \neq i \\ 0 & : k = i \end{cases} \quad (15)$$

Proposition 1 suggests an efficient way of obtaining the solution to (14). According to (15),  $\alpha^{i(t)} = \mathbf{v}^*$  can be obtained by searching over a candidate set of size  $K + 1$ , where  $K$  is the number of nearest neighbors to construct the KNN graph  $\mathbf{S}$  for regularized  $\ell^0$ -graph.

Similar to Theorem 2, the sequence  $\{F(\alpha^{i(t)})\}_t$  is decreasing. The iterative proximal method starts from  $t = 1$  and continue until the sequence  $\{F(\alpha^{i(t)})\}_t$  converges or

maximum iteration number is achieved. When the proximal method converges or terminates for each  $\alpha^i$ , the step of coordinate descent for  $\alpha^i$  is finished and the optimization algorithm proceeds to optimize other sparse codes. Each iteration of coordinate descent solves (12) for  $i = 1 \dots n$  sequentially, and it terminates when maximum iteration number is reached or converges under some stopping threshold on the change of the objective function (11).

## 5. Experimental Results

The superior clustering performance of  $\ell^0$ -graph is demonstrated in this section with extensive experimental results, and we also show the effectiveness of regularized  $\ell^0$ -graph. We compare our  $\ell^0$ -graph to K-means (KM), Spectral Clustering (SC),  $\ell^1$ -graph, Sparse Manifold Clustering and Embedding (SMCE) [7]. Moreover, we derive the OMP-graph, which builds the sparse graph in the same way as  $\ell^0$ -graph except that it solves the following optimization problem by Orthogonal Matching Pursuit (OMP) to obtain the sparse code:

$$\min_{\alpha^i} \|\mathbf{x}_i - \mathbf{X}\alpha^i\|_F^2 \quad s.t. \quad \|\alpha^i\|_0 \leq T, \alpha_i^i = 0, i = 1, \dots, n \quad (16)$$

$\ell^0$ -graph is also compared to OMP-graph to show the advantage of the proposed proximal method in the previous sections. By adjusting the parameters,  $\ell^1$ -graph and SSC solve the same problem and generate equivalent results, so we report their performance under the same name “ $\ell^1$ -graph”.

### 5.1. Evaluation Metric

Two measures are used to evaluate the performance of the clustering methods, i.e. the accuracy and the Normalized Mutual Information (NMI) [31]. Let the predicted label of the datum  $\mathbf{x}_i$  be  $\hat{y}_i$  which is produced by the clustering method, and  $y_i$  is its ground truth label. The accuracy is defined as

$$\text{Accuracy} = \frac{\mathbb{I}_{\Omega(\hat{y}_i) \neq y_i}}{n} \quad (17)$$

where  $\mathbb{I}$  is the indicator function, and  $\Omega$  is the best permutation mapping function by the Kuhn-Munkres algorithm [21]. The more predicted labels match the ground truth ones, the more accuracy value is obtained.

Let  $\hat{X}$  be the index set obtained from the predicted labels  $\{\hat{y}_i\}_{i=1}^n$  and  $X$  be the index set from the ground truth labels  $\{y_i\}_{i=1}^n$ . The mutual information between  $\hat{X}$  and  $X$  is

$$MI(\hat{X}, X) = \sum_{\hat{x} \in \hat{X}, x \in X} p(\hat{x}, x) \log_2 \left( \frac{p(\hat{x}, x)}{p(\hat{x})p(x)} \right) \quad (18)$$

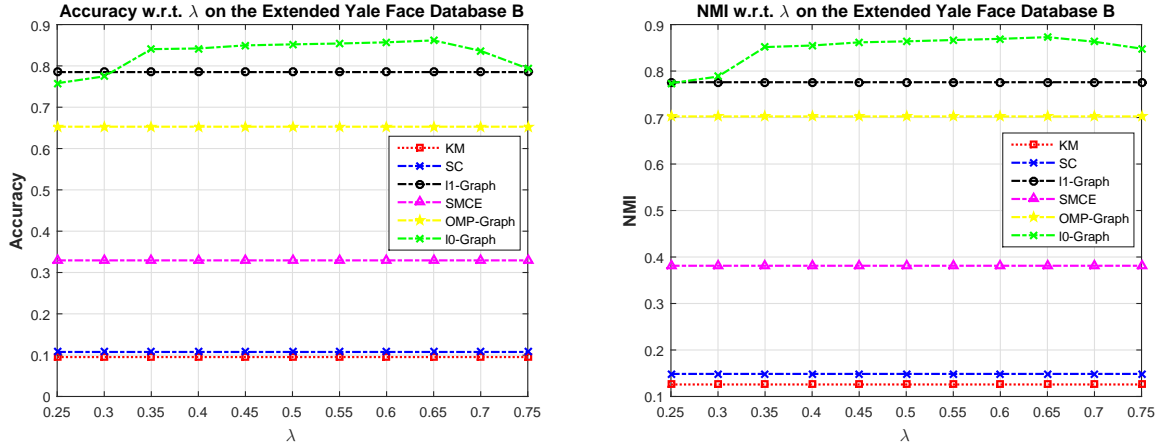
where  $p(\hat{x})$  and  $p(x)$  are the margined distribution of  $\hat{X}$  and  $X$  respectively, induced from the joint distribution  $p(\hat{x}, x)$

Table 2. Clustering Results on Ionosphere and MNIST Handwritten Digits Database

Data Set	Measure	KM	SC	$\ell^1$ -Graph	SMCE	OMP-Graph	$\ell^0$ -Graph
Ionosphere	AC	0.7097	0.7350	0.5128	0.6809	0.6353	<b>0.7692</b>
	NMI	0.1287	0.2155	0.1165	0.0871	0.0299	<b>0.2609</b>
MNIST	AC	0.5621	0.4922	0.4948	0.5784	0.5754	<b>0.6590</b>
	NMI	0.5113	0.4755	0.5210	0.6332	0.5463	<b>0.6709</b>

Table 3. Clustering Results on COIL-20 Database.  $c$  in the left column is the cluster number, i.e. the first  $c$  clusters of the entire data are used for clustering.  $c$  has the same meaning in Table 4 and Table 5.

COIL-20 # Clusters	Measure	KM	SC	$\ell^1$ -Graph	SMCE	OMP-Graph	$\ell^0$ -Graph
$c = 4$	AC	0.6632	0.6701	1.0000	0.7639	0.9271	<b>1.0000</b>
	NMI	0.5106	0.5455	1.0000	0.6741	0.8397	<b>1.0000</b>
$c = 8$	AC	0.5130	0.4462	0.7986	0.5365	0.6753	<b>0.9705</b>
	NMI	0.5354	0.4947	0.8950	0.6786	0.7656	<b>0.9638</b>
$c = 12$	AC	0.5885	0.4965	0.7697	0.6806	0.5475	<b>0.8310</b>
	NMI	0.6707	0.6096	0.8960	0.8066	0.6316	<b>0.9149</b>
$c = 16$	AC	0.6579	0.4271	0.8273	0.7622	0.3481	<b>0.9002</b>
	NMI	0.7555	0.6031	0.9301	0.8730	0.4520	<b>0.9552</b>
$c = 20$	AC	0.6554	0.4278	0.7854	0.7549	0.3389	<b>0.8472</b>
	NMI	0.7630	0.6217	0.9148	0.8754	0.4853	<b>0.9428</b>

Figure 1. Clustering performance with different values of  $\lambda$ , i.e. the weight for the  $\ell^0$ -norm, on the Extended Yale Face Database B. Left: Accuracy; Right: NMI

over  $\hat{X}$  and  $X$ . Let  $H(\hat{X})$  and  $H(X)$  be the entropy of  $\hat{X}$  and  $X$ , then the normalized mutual information (NMI) is defined as below:

$$NMI(\hat{X}, X) = \frac{MI(\hat{X}, X)}{\max\{H(\hat{X}), H(X)\}} \quad (19)$$

It can be verified that the normalized mutual information takes values in  $[0, 1]$ . The accuracy and the normalized mutual information have been widely used for evaluating the performance of the clustering methods [30, 4, 31].

## 5.2. Clustering on UCI Data Set and MNIST Handwritten Digits Database

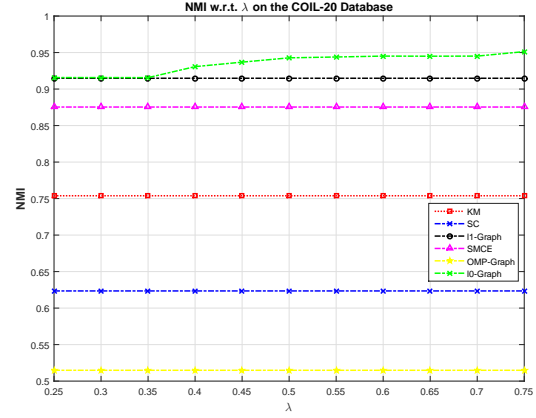
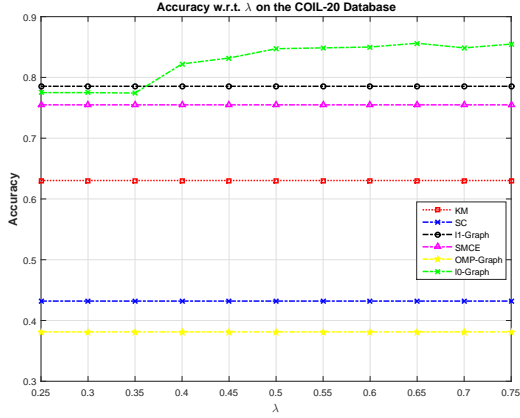
In this subsection, we conduct experiments on the Ionosphere data from UCI machine learning repository [1] and the MNIST database of handwritten digits. The information of these two data sets are in Table 6. MNIST handwritten digits database has a total number of 70000 samples for digits from 0 to 9. The digits are normalized and centered in a fixed-size image. For MNIST data set, we randomly select 500 samples for each digit to obtain a subset of MNIST data consisting of 5000 samples. The random sampling is performed for 10 times and the average clustering performance is recorded. The clustering results on the two data

Table 4. Clustering Results on COIL-100 Database.

COIL-100 # Clusters	Measure	KM	SC	$\ell^1$ -Graph	SMCE	OMP-Graph	$\ell^0$ -Graph
c = 20	AC	0.5850	0.4514	0.5757	0.6208	0.4243	<b>0.9264</b>
	NMI	0.7456	0.6700	0.7980	0.7993	0.5258	<b>0.9681</b>
c = 40	AC	0.5791	0.4139	0.5934	0.6038	0.2340	<b>0.8472</b>
	NMI	0.7691	0.6681	0.7962	0.7918	0.4378	<b>0.9471</b>
c = 60	AC	0.5371	0.3389	0.5657	0.5887	0.1905	<b>0.8326</b>
	NMI	0.7622	0.6343	0.8162	0.7973	0.3690	<b>0.9352</b>
c = 80	AC	0.5048	0.3115	0.5271	0.5835	0.2247	<b>0.7899</b>
	NMI	0.7474	0.6088	0.8006	0.8006	0.4173	<b>0.9218</b>
c = 100	AC	0.4996	0.2835	0.5275	0.5639	0.1667	<b>0.7683</b>
	NMI	0.7539	0.5923	0.8041	0.8064	0.3757	<b>0.9182</b>

Table 5. Clustering Results on the Extended Yale Face Database B.

Yale-B # Clusters	Measure	KM	SC	$\ell^1$ -Graph	SMCE	OMP-Graph	$\ell^0$ -Graph
c = 10	AC	0.1782	0.1922	0.7580	0.3672	0.7375	<b>0.8406</b>
	NMI	0.0897	0.1310	0.7380	0.3266	0.7468	<b>0.7695</b>
c = 15	AC	0.1554	0.1706	0.7620	0.3761	0.7532	<b>0.7987</b>
	NMI	0.1083	0.1390	0.7590	0.3593	0.7943	<b>0.8183</b>
c = 20	AC	0.1200	0.1466	0.7930	0.3526	0.7813	<b>0.8273</b>
	NMI	0.0872	0.1183	0.7860	0.3771	0.8172	<b>0.8429</b>
c = 30	AC	0.1096	0.1209	0.8210	0.3470	0.7156	<b>0.8633</b>
	NMI	0.1159	0.1338	0.8030	0.3927	0.7260	<b>0.8762</b>
c = 38	AC	0.0954	0.1077	0.7850	0.3293	0.6529	<b>0.8480</b>
	NMI	0.1258	0.1485	0.7760	0.3812	0.7024	<b>0.8612</b>

Figure 2. Clustering performance with different values of  $\lambda$ , i.e. the weight for the  $\ell^0$ -norm, on the COIL-20 Database. Left: Accuracy; Right: NMI

sets are shown in Table 2.

Table 6. Two UCI data sets and MNIST Handwritten Digits Database in the experiments

	Heart	Ionosphere	MNIST
# of instances	270	351	70000
Dimension	13	34	1024
# of classes	2	2	10

### 5.3. Clustering On COIL-20 and COIL-100 Database

COIL-20 Database has 1440 images of 20 objects in which the background has been removed, and the size of each image is  $32 \times 32$ , so the dimension of this data is 1024. COIL-100 Database contains 100 objects with 72 images of size  $32 \times 32$  for each object. The images of each object were taken 5 degrees apart when the object was rotated on

a turntable. The clustering results on these two data sets are shown in Table 3 and Table 4 respectively. We observe that  $\ell^0$ -graph performs consistently better than all other competing methods. On COIL-100 Database, SMCE renders slightly better results than  $\ell^1$ -graph on the entire data due to its capability of modeling non-linear manifolds.

#### 5.4. Clustering On Extended Yale Face Database B

The Extended Yale Face Database B contains face images for 38 subjects with 64 frontal face images taken under different illuminations for each subject. The clustering results are shown in Table 5. We can see that  $\ell^0$ -graph achieves significantly better clustering result than  $\ell^1$ -graph, which is the second best method on this data.

#### 5.5. Improved $\ell^0$ -Graph with Regularization

In this subsection, we investigate the performance of regularized  $\ell^0$ -graph. We empirically set  $\mathbf{S}$  to be the adjacency matrix of 5-NN graph and  $\gamma = 0.1$  as the default parameter setting for regularized  $\ell^0$ -graph in (11). We conduct comparison experiments on the UCI Heart data whose information is in Table 2, the Extended Yale Face Database B and UMIST Face Database. The UMIST Face Database consists of 575 images of size  $112 \times 92$  for 20 people. Each person is shown in a range of poses from profile to frontal views. The clustering results are shown in Table 7. The better results of regularized  $\ell^0$ -graph are due to the fact that it promotes common neighbors for nearby data so as to produce a more aligned similarity graph and alleviate the graph connectivity issue.

Table 7. Clustering Performance of Regularized  $\ell^0$ -Graph

Data Set	Measure	$\ell^0$ -Graph	R $\ell^0$ -Graph
Heart	AC	0.5111	<b>0.6444</b>
	NMI	0.0064	<b>0.0590</b>
Extended Yale B	AC	0.8480	<b>0.8521</b>
	NMI	0.8612	<b>0.8634</b>
UMIST Face	AC	0.6730	<b>0.7078</b>
	NMI	0.7924	<b>0.8153</b>

#### 5.6. Parameter Setting

We use the sparse codes generated by  $\ell^1$ -graph with the weighting parameter  $\lambda_{\ell^1} = 0.1$  in (3), which is the default value suggested in [8], to initialize  $\ell^0$ -graph, and set  $\lambda = 0.5$  for  $\ell^0$ -graph empirically throughout all the experiments in this section. We observe that the average number of non-zero elements of the sparse code for each data point is around 3 for most data sets. The maximum iteration number  $M = 100$  and the stopping threshold  $\varepsilon = 10^{-6}$ . For OMP-graph, we tune the parameter  $T$  in (16) to control the sparsity of the generated sparse codes such that the aforementioned average number of non-zero elements of

the sparse code matches that of  $\ell^0$ -graph. For  $\ell^1$ -graph, the weighting parameter for the  $\ell^1$ -norm is chosen from  $[0.1, 1]$  for the best performance.

We investigate how the clustering performance on the Extended Yale Face Database B and COIL-20 Database changes by varying the weighting parameter  $\lambda$  for  $\ell^0$ -graph, and illustrate the result in Figure 1 and Figure 2 respectively. We observe that the performance of  $\ell^0$ -graph is much better than other algorithms over a relatively large range of  $\lambda$ , revealing the robustness of our algorithm with respect to the weighting parameter  $\lambda$ .

#### 5.7. Efficient Parallel Computing by CUDA Implementation

We have implemented  $\ell^0$ -graph, regularized  $\ell^0$ -graph in CUDA C programming language on NVIDIA K40. Both the MATLAB and CUDA implementation will be available for downloading. We compare the running time of  $\ell^0$ -graph in MATLAB implementation and CUDA C implementation on the Extended Yale Face Database B data, on a workstation with 2 Intel Xeon X5650 2.67 GHz CPU, 48 GB memory and one NVIDIA K40 graphics card. MATLAB implementation takes 48.51 seconds while the CUDA implementation only takes 1.68 seconds, with a speedup of 28.87 times.

Due to the limited space, we have put additional experimental results in the supplementary document for this paper, such as the application of  $\ell^0$ -graph on semi-supervised learning, and the parameter sensitivity for regularized  $\ell^0$ -graph.

### 6. Conclusion

We propose a novel  $\ell^0$ -graph for data clustering in this paper. In contrast to the existing sparse subspace clustering method such as Sparse Subspace Clustering and  $\ell^1$ -graph,  $\ell^0$ -graph features  $\ell^0$ -induced almost surely subspace-sparse representation under milder assumptions on the subspaces and random data generation. The objective function of  $\ell^0$ -graph is optimized using a proposed proximal method. Convergence of this proximal method is proved, and extensive experimental results on various real data sets demonstrate the effectiveness and superiority of  $\ell^0$ -graph over other competing methods. To improve the graph connectivity, we propose regularized  $\ell^0$ -graph whose effectiveness is also demonstrated on real data sets.

### References

- [1] D. N. A. Asuncion. UCI machine learning repository, 2007.
- [2] C. Bao, H. Ji, Y. Quan, and Z. Shen. L0 norm based dictionary learning by proximal methods with global convergence. In *2014 IEEE Conference on Computer Vision and Pattern*



- Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014, pages 3858–3865, 2014. [3](#)
- [3] J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Math. Program.*, 146(1-2):459–494, Aug. 2014. [3](#), [4](#)
- [4] B. Cheng, J. Yang, S. Yan, Y. Fu, and T. S. Huang. Learning with l1-graph for image analysis. *IEEE Transactions on Image Processing*, 19(4):858–866, 2010. [1](#), [2](#), [6](#)
- [5] H. Cheng, Z. Liu, L. Yang, and X. Chen. Sparse representation and learning in visual recognition: Theory and applications. *Signal Process.*, 93(6):1408–1425, June 2013. [2](#)
- [6] E. L. Dyer, A. C. Sankaranarayanan, and R. G. Baraniuk. Greedy feature selection for subspace clustering. *Journal of Machine Learning Research*, 14:2487–2517, 2013. [2](#)
- [7] E. Elhamifar and R. Vidal. Sparse manifold clustering and embedding. In *NIPS*, pages 55–63, 2011. [1](#), [5](#)
- [8] E. Elhamifar and R. Vidal. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(11):2765–2781, 2013. [1](#), [2](#), [4](#), [8](#)
- [9] C. Fraley and A. E. Raftery. Model-Based Clustering, Discriminant Analysis, and Density Estimation. *Journal of the American Statistical Association*, 97(458):611–631, June 2002. [1](#)
- [10] M. Hyder and K. Mahata. An approximate l0 norm minimization algorithm for compressed sensing. In *Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on*, pages 3365–3368, April 2009. [3](#)
- [11] R. Jenatton, J. Mairal, F. R. Bach, and G. R. Obozinski. Proximal methods for sparse hierarchical dictionary learning. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 487–494, 2010. [1](#)
- [12] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma. Robust recovery of subspace structures by low-rank representation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(1):171–184, Jan. 2013. [1](#), [4](#)
- [13] G. Liu, Z. Lin, and Y. Yu. Robust subspace segmentation by low-rank representation. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel*, pages 663–670, 2010. [1](#), [4](#)
- [14] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *J. Mach. Learn. Res.*, 11:19–60, Mar. 2010. [1](#)
- [15] J. Mairal, F. R. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Supervised dictionary learning. In *Advances in Neural Information Processing Systems 21, Proceedings of the Twenty-Second Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 8-11, 2008*, pages 1033–1040, 2008. [1](#)
- [16] L. Mancera and J. Portilla. L0-norm-based sparse representation through alternate projections. In *Image Processing, 2006 IEEE International Conference on*, pages 2089–2092, Oct 2006. [2](#)
- [17] B. Nasihatkon and R. Hartley. Graph connectivity in sparse subspace clustering. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 2137–2144, June 2011. [4](#)
- [18] A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In *NIPS*, pages 849–856, 2001. [1](#), [2](#)
- [19] D. Park, C. Caramanis, and S. Sanghavi. Greedy subspace clustering. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2753–2761, 2014. [1](#), [2](#), [3](#), [4](#)
- [20] X. Peng, Z. Yi, and H. Tang. Robust subspace clustering via thresholding ridge regression. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 3827–3833. AAAI, 2015. [1](#)
- [21] D. Plummer and L. Lovász. *Matching Theory*. North-Holland Mathematics Studies. Elsevier Science, 1986. [5](#)
- [22] M. Soltanolkotabi and E. J. Cands. A geometric analysis of subspace clustering with outliers. *Ann. Statist.*, 40(4):2195–2238, 08 2012. [1](#), [3](#), [4](#)
- [23] J. A. Tropp. Greed is good: algorithmic results for sparse approximation. *IEEE Transactions on Information Theory*, 50(10):2231–2242, 2004. [2](#), [3](#)
- [24] R. Vidal. Subspace clustering. *Signal Processing Magazine, IEEE*, 28(2):52–68, March 2011. [1](#)
- [25] Y.-X. Wang, H. Xu, and C. Leng. Provable subspace clustering: When lrr meets ssc. In C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 64–72. Curran Associates, Inc., 2013. [1](#), [2](#), [3](#), [4](#)
- [26] S. Yan and H. Wang. Semi-supervised learning by sparse representation. In *SDM*, pages 792–801, 2009. [1](#), [2](#)
- [27] J. Yang, K. Yu, Y. Gong, and T. S. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *CVPR*, pages 1794–1801, 2009. [2](#)
- [28] Y. Yang, Z. Wang, J. Yang, J. Han, and T. Huang. Regularized l1-graph for data clustering. In *Proceedings of the British Machine Vision Conference*. BMVA Press, 2014. [1](#), [2](#), [5](#)
- [29] T. Zhang, B. Ghanem, S. Liu, C. Xu, and N. Ahuja. Low-rank sparse coding for image classification. In *IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013*, pages 281–288, 2013. [2](#)
- [30] M. Zheng, J. Bu, C. Chen, C. Wang, L. Zhang, G. Qiu, and D. Cai. Graph regularized sparse coding for image representation. *IEEE Transactions on Image Processing*, 20(5):1327–1336, 2011. [5](#), [6](#)
- [31] X. Zheng, D. Cai, X. He, W.-Y. Ma, and X. Lin. Locality preserving clustering for image database. In *Proceedings of the 12th Annual ACM International Conference on Multimedia, MULTIMEDIA '04*, pages 885–891, New York, NY, USA, 2004. ACM. [5](#), [6](#)